

Data Mining & Warehousing :-

Data Mining :-

→ Knowledge mining from large amount of data.

→ Knowledge discovery from Data (KDD)

Extracting the knowledge from the large amount of data.

knowledge discovery process is an iterative sequence of the following steps :-

- Step 1 → ① Data cleaning
- ② Data Integration (maybe data present at diff' places so you have to perform integration)
- ③ Data Selection
- ④ Data Transformation. → means you have to select relevant data.
- ⑤ Data Mining
- ⑥ Pattern evaluation
- ⑦ Knowledge presentation

apply diff' algo.
to perform
data mining.

in this you
present knowledge
in visualization
form e.g. graph,
chart.

- Temporal Database :-
- Sequence Database
- Time-Series "

(1) Temporal database

↳ stores relational data that includes the time related attributes.

↳ timestamp.

② Sequence database:
It records sequence of events or activity happens

③ Time-series database:-

↳ stores the sequence of values or events over repeated measurement of time (minute/second/mill sec./daily/weekly etc.)

e.g:-

Diff' b/w Temporal & Time-Series database:-

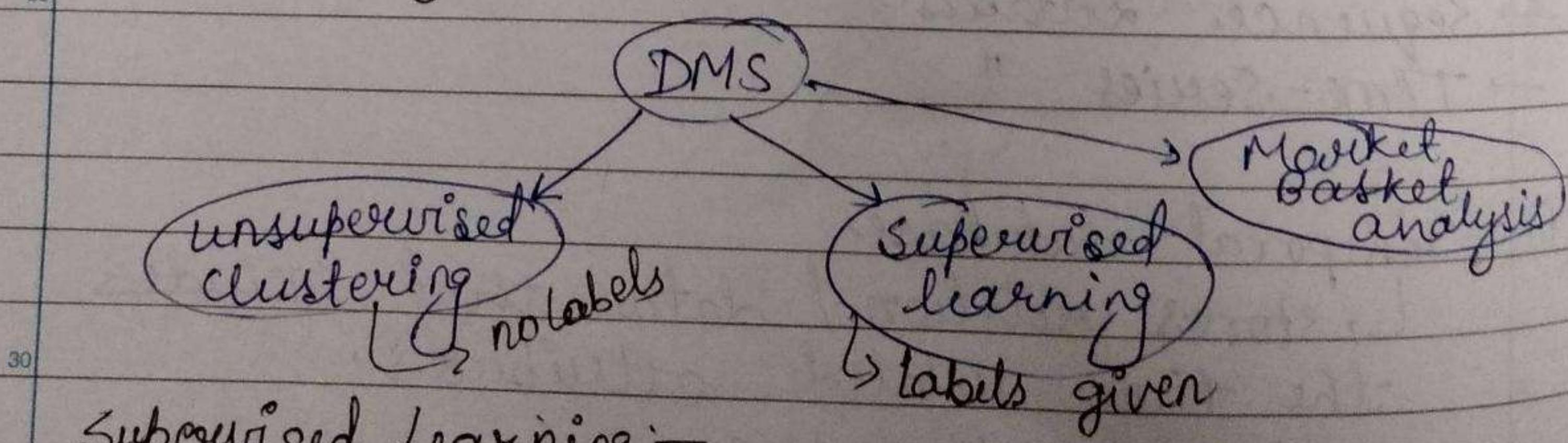
time interval is
not fixed

time interval is
fixed

→ Data Mining

data mining (KDD) is the process of analysing data from different perspective & summarising it into useful information that can be used to → increase revenue.
→ cut cost or both.

→ Data Mining Strategies:-



Supervised Learning:-

→ Mind. choose what it believe to be the

Page :
Date :

defining concept, feature & form our own
classification model.

→ Supervision

 + → Labels

 ↳ given to instances.

→ induction based learning

Supervised Learning

Classification

Regression

Classification:

 + → the output is having certain defined labels
(discrete values).

 + → Binary or multiclass. [O/P: → yes/no
 + → Red/green/blue].

Regression

 + → here the output is having continuous
values. e.g. temp, pressure, ~~wind dir~~, humidity
 wind speed

 → Predict the sale,

 → Price of house.

 " " star

→ Unsupervised Learning.

 + → it builds model from data without pre-defined classes.

 + → data instances are grouped together based
on the similarity scheme defined by
clustering system.

→ clustering

↳ grouping the data into clusters, so that objects within a cluster have high similarity.

→ but are very dissimilar to objects in other clusters.

→ dissimilarity, dissimilarity is accessed based on the attribute values, that are describing the object.

→ Distance measure, is generally used.

→ Supervised]⇒ semi supervised Learning
Unsupervised

↳ majority of data not having labels,
only small set of data having labels

→ First we apply clustering algo
↳ based on clustering labelled it
↓
labelled data

→ Market Basket Analysis :-

↳ determines which sets of products tend to be purchased together.

Mobile → Mobile cover | e.g.

Mobile → Tempered glass.

+ → Rules
+ → Sales

→ Association Rule.

↳ Recommendations

A → B
Item Item

→ Take few numbers

13, 18, 13, 14, 13, 16, 14, 21, 13

$$\text{Mean} = \frac{\text{Sum of numbers}}{\text{Total no.}} = \frac{135}{9} = 15$$

Median

13, 13, 13, 13, 13, 14, 14, 16, 18, 21

Find the median = 14.

it is the middle value of the dataset when it has been arranged in order.

e.g. ② 12, 18, 16, 21, 10, 13, 17, 19

10, 12, 13, 17, 18, 19, 21

Median = 16.5

Mean = 15.75

Median = 16.5

Camlin

if the dataset is even
median is the average of two middle values

→ Mode :-

48, 44, 48, 45, 42, 49, 48

mode = 48

data values in the dataset that occurs most often.

e.g:- dataset is given like this

age	frequency
1 - 5	200
5 - 15	450
15 - 20	300
20 - 50	1500
50 - 80	700
80 - 110	44

Here median comes when
($\frac{N}{2}$ which is 1597)

calculate the approx median value of the data.

Sol? median = $L_1 + \left(\frac{N/2 - (\sum freq)_l}{freq_{median}} \right) width$

where, $L_1 \rightarrow$ lower boundary of the median interval.
 $N \rightarrow$ is the no. of values in the dataset

$(\sum freq)_l \rightarrow$ sum of frequencies of all the intervals that are lower than the median interval

$freq_{median} \rightarrow$ is the freq. of the median interval

width \rightarrow It is the width of the median interval.

10,000, median = 1500

median = 152.5

N=3194, median = 1507

$L_1 = 20, (\geq f_{req})_1 = 950, f_{req, median} = 1500$

width = 30

median = $20 + \frac{1}{2}$

* Data Preprocessing :-

→ Data redundancy

↳ when same data present at more than one places.

→ Data inconsistency

↳ if data is present at multiple places then if value of same data is diff. at diff. places so it is inconsistency.

→ (1) Data cleaning

+ incomplete, noisy, inconsistent

→ fill in the missing values

→ smooth out the noise while identifying outliers

→ correct inconsistency in the data.

Missing values

→ handle it

① → ignore the tuple

Canon

F_1	F_2	F_3	F_4	F_5	Class Label
Adult	Gender	F_3	F_4	F_5	X

if tuple is missing ignore the tuple.

F_1	F_2	F_3	F_4	F_5	Class Label
1	1	-	X	✓	[not ignore]

F_1	F_2	F_3	F_4	F_5	Class Label
X	X	X	X	✓	✓

when most of the entries in tuple is missing so ignore the tuple.

- ② Fill in the missing values manually:-
[when only few values are missing]

- ③ Use a global constant to fill in the missing value.
e.g. unknown, 0, ∞ , etc., not completely good.

- ④ Use the attribute mean to fill in the missing value.

F_1	Income		class
34	100		
26	200		
33	Mean		
28	300		
33	150		
11	200		

- ⑤ Use the attribute mean for all samples belonging to the same class as the given tuple.

F_1	Income	f_n	class
34	100	↑	A
26	200	↑	A
31	300	↑	A
29	meanA	↓	A
30	400	↑	B
11	50	↓	B
12	meanB	↓	B
19	20		B

⑥ Use the most probable values to fill in the missing values.
e.g.: Regression, Bayesian, decision tree

* Standard Deviation :-

1, 2, 3, 4, 5, 6.

$$\text{S.D} \quad \sigma = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}}$$

Mean = 3.5, SD = 1.7

→ Data Cleaning

- ↳ Missing Values
- ↳ Noisy Data

→ Noisy Data :-

- ↳ random error
- ↳ Binning
- ↳ Regression
- ↳ Clustering

Binning :-

- ↳ It smooths data values by consulting its neighbourhood. i.e. the value around it.

e.g. 4, 8, 15, 21, 21, 24, 25, 28, 34

- Put the data values into buckets/bins.

→ Local smoothing

- ↳ smoothing by bin means
- ↳ smoothing by bin boundaries

Bin 1: 4, 8, 15

Bin 2: 21, 21, 24

Bin 3: 25, 28, 34

① Smoothing by bin means:-

Bin 1: 9, 9, 9

Bin 2: 22, 22, 22

Bin 3: 29, 29, 29.

② Smoothing by bin boundaries:-

Bin 1: 4, 4, 15 , min=4, max=15

Bin 2: 21, 21, 24 , min=21, max=24

Bin 3: 25, 25, 34 , min=25, max=34

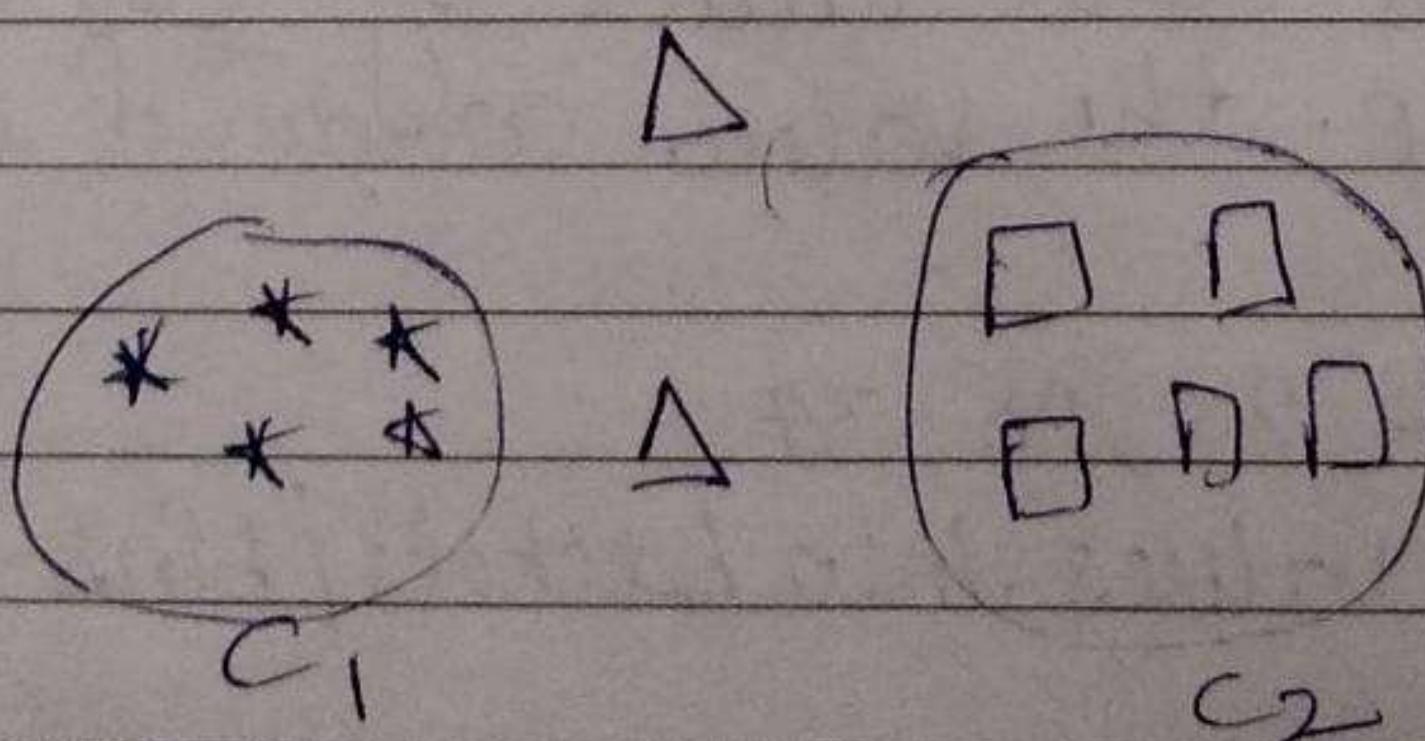
Min^m & Max^m values are identified \Rightarrow bin boundaries

Each bin value is then replaced by the closest boundary value.

* Regression

→ data can be smoothed by fitting it into a function.

clustering:



Outliers

→ They fall outside the set of clusters.

→ Data Cleaning :-

- discrepancy detection
- metadata

→ Data scrubbing tools

use domain knowledge to detect & correct data.

→ Data Auditing tool

It finds rules & relationship

Data Mining Techniques

→ Correlation Analysis :-

relationship b/w numerical attributes

↓

Pearson correlation.

$$\gamma_{A,B} = \frac{\sum_{i=1}^N (a_i - \bar{A})(b_i - \bar{B})}{\sqrt{\sum_{i=1}^N (a_i - \bar{A})^2} \sqrt{\sum_{i=1}^N (b_i - \bar{B})^2}}$$

N → no. of tuples

a_i & b_i are respective values of the tuple.

$\gamma_{A,B} > 0$, positively correlated i.e. if A ↑ ^{then} B ↑

$$-1 \leq \gamma_{A,B} \leq +1$$

~~Ex~~ → -ve value \Rightarrow Negatively correlated.

+ve value \Rightarrow Positively "

0. \rightarrow A, B are independent.
(zero) \wedge no correlation.

Camlin

to write

$$\bar{A} = 40.12$$

$$\bar{B} = 9.12$$

Page :
Date :

Tree Height	Trunk diameter
48	8
35	9
49	7
27	6
33	13
21	7
45	11
51	12

Find the correlation b/w Tree Height & Trunk diameter.

$$\text{Ans: } \approx 0.89$$

+ positively correlated.

U solve using $r_{AB} =$

** Data transformation:-

⇒ transformed or consolidated into one form appropriate for mining.

① → Smoothing

- ↳ Removal of noise
 - ↳ Binning
 - ↳ Regression
 - ↳ Clustering

② → Aggregation

- ↳ Aggregation oppn?
 - ↳ daily → weekly

③ → Generalization

- We are replacing the low level data by the higher level concept through the concept hierarchy.

Ex :- City can → state

① Normalization :-

→ Data are scaled so as to fall within in a specified range.

-1 to +1

OR

0 to 1.

② Attribute construction

→ new attributes are constructed from the given set of attributes.

** Normalization Techniques

→ Min-max normalization

$\min_A \rightarrow \text{min}^m \text{ value of an attribute } A.$

$\max_A \rightarrow \text{max}^m \quad " \quad " \quad " \quad A.$

↓

In the new range,

$\text{new-min}_A \rightarrow \text{new min. val of an attribute } A$

$\text{new-max}_A \rightarrow \text{new max value of an attribute } A.$

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new max}_A - \text{new-min}_A) + \text{new-min}_A$$

e.g. min & max value for attribute income is

₹ 12000/- & ₹ 98000/- . Map into range [0.0 to 1.0].

& find the value of ₹ 73600/-

sol.

$$= \frac{73600 - 12000}{98000 - 12000} * (1 - 0) + 0$$

$$= 0.716$$

Camlin

to find

→ Z-score normalization: or zero-mean normalization:-

Value of an attribute A are normalized

↓

based on mean & standard deviation of A.
 initial value $\rightarrow v$
 After normalization $\rightarrow v'$

$$v' = \frac{v - \bar{A}}{\sigma_A}$$

$\bar{A} \rightarrow$ mean

$\sigma_A \rightarrow$ Standard deviation.

e.g. attribute income

mean $\rightarrow ₹ 54000/-$

SD $\rightarrow ₹ 6000/-$

with z-score norm find the value of 73600

$$v' = \frac{73600 - 54000}{60000} = 0.32666$$

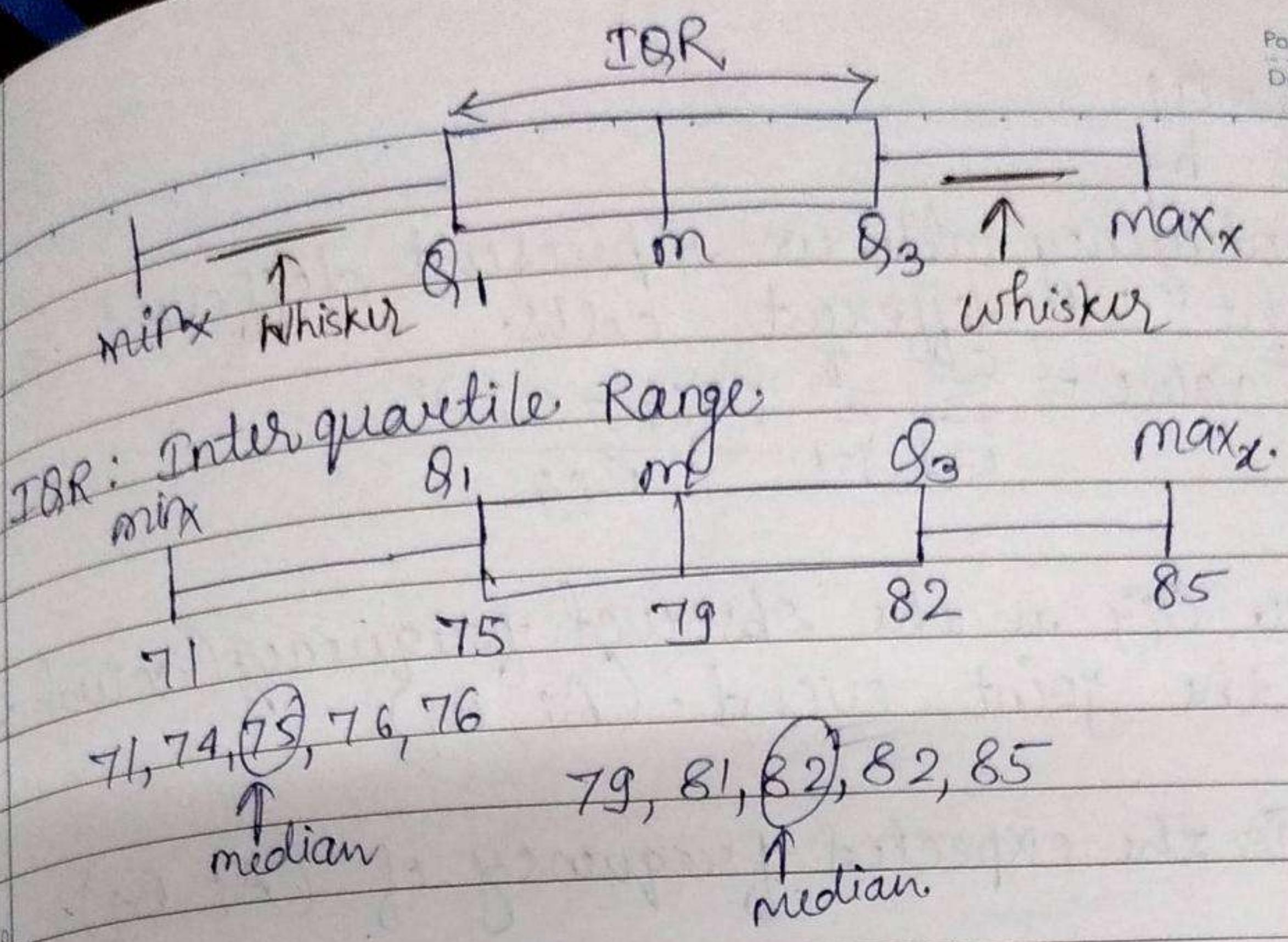
* Boxplot :-

76, 79, 76, 74, 75, 71, 85, 82, 82, 79, 81

→ Arrange the value in the ascending order.

71, 74, 75, 76, 76, 79, 79, 81, 82, 82, 85

→ It gives 5 number summary of the data distribution.



→ Pearson Correlations :-

$$-1 \leq r_{AB} \leq 1.$$

χ^2 (Chi Square Test)

→ It finds the correlation b/w two categorical attributes.

Let's say those attributes are A & B.

→ A → 'c' distinct values

$$a_1, a_2, a_3, \dots, a_c$$

→ B → 's' distinct values

$$b_1, b_2, b_3, \dots, b_s$$

represented in the form of 'contingency table'.

Contingency Table

'c' values of 'A' making up the col. & 's' value of 'B' making up the rows.

$(A_i, B_j) \rightarrow$ denote an event that

$$A \rightarrow a_i$$

$$B \rightarrow b_j$$

Camlin

to find
out
the
value

$A = a_i$ $B = b_j$

→ In contingency table we represent different joint event in different cells.

$$\chi^2 \text{ value} = \sum_{i=1}^{C} \sum_{j=1}^{R} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

where, O_{ij} is the observed frequency (actual count) of the joint event (A_i, B_j) .

E_{ij} is the expected frequency of (A_i, B_j) .

$$E_{ij} = \frac{\text{Count}(A=a_i) * \text{Count}(B=b_j)}{N}$$

where $N \rightarrow$ no. of data-tuples

$\text{Count}(A=a_i) \rightarrow$ No. of tuples having value a_i for A.

$\text{Count}(B=b_j) \rightarrow$ No. of tuples having value b_j for B.

e.g.		male	female	total
	fiction	250	200	450
observed frequency	Non-fiction	50	1000	1050
		300	1200	1500

$$e_{11} = \frac{\text{count(male)} * \text{count(fiction)}}{N}$$

$$= \frac{300 * 450}{1500} = 90$$

$$e_{12} = \frac{1200 * 450}{1500} = 360$$

$$e_{21} = \frac{1200 * 1050}{1500} = 840$$

$$e(\text{male, nonfiction}) = \frac{2}{\frac{300 \times 1050}{1500}} = 210$$

$$\begin{aligned} \chi^2 &= \frac{(250-90)^2}{90} + \frac{(50-210)^2}{210} + \frac{(200-360)^2}{360} \\ &\quad + \frac{(1000-840)^2}{840} \\ &= 507 \\ &\approx 508 \end{aligned}$$

Hypothesis \rightarrow There is no relationship b/w categorical variable A & categorical variable B.

$H_1 \rightarrow$ There is some relationship b/w these two categorical variable.

Step ①

Find the ~~etg~~ degree of freedom

$$\hookrightarrow (r-1) \times (c-1)$$

$$dof = (2-1) \times (2-1) = 1$$

Step ②

Significance level \rightarrow always given

1% or 5%

χ^2 , dof, SL
Value

- \rightarrow If the value in the ~~etg~~ chart is less than calculated χ^2 then we accept H_1 .
- \rightarrow If the calculated χ^2 is less than we accept H_0 .

class

Data warehouse:-

- OLTP
- OLAP systems

OLTP :- (Online Transaction processing)
↳ we are using in day to day operation.

OLAP: (Online Analytical processing) :-

We put historical data after integration / transformation and various other operations.

→ OLAP → Analysis / Decision Making

↳ Managers / Owners
or other Policy makers.

→ Data warehouse:-

- to systematically organize
- Understand
- Use their data to make strategic decisions

According to William H. Inman :- (Characteristics of data warehouse)

→ Subject oriented

→ Integrated

→ Time-variant &

→ Non-volatile.

in support of management decision making process

→ Subject oriented:-

↳ Centered around major subjects etc
ex :- sales, suppliers.

Integrated :-

→ Integrating data from multiple sources.
These sources may be homogeneous or heterogeneous.

Time-Variants

→ Data is stored from the historical perspective.

→ Every key structure in DW contains either implicitly or explicitly an element of time.

Non-volatile :-

→ We perform two operations.

 → Initial loading

 → Access of data.

→ Major difference b/w OLTP & OLAP:-

OLTP → day to day operation.

OLAP → Used by decision makers.

Data contents :-

OLTP → current data, detailed

OLAP → large amount of historical, aggregated data.

Database Design -

OLTP → ER models

OLAP → star or snowflake.

fact-constellation.

Access Pattern:

OLTP → short, atomic transactions

OLAP → read only operations.

Data Cubes :-

↳ multidimensional data models
↳ Data cubes.

- A data cube allows data to be modeled and viewed in multiple dimensions.
 → dimensions
 → facts
 → entities over which an organization wants to keep records.

For sales subject
dim → time, items, branch, location etc.

- Each dimⁿ has table associated with it.
 ↓
 dimension table:

- The subject around which multidimⁿ data model is organized is represented by FACT table.

EX:- 2D- view of sales data

dim → time & item

Sales from branches in Patna.

location = "Patna"

item type

Quarter	time	Home enter	computer	laptop	mobile
Q1		609	961	106	-
Q2		514	102	-	-
Q3		623	103	107	-
Q4		691	104	-	-

→ Can organize the data in 3D.

Page :
Date :

location = "Patna"

items
[H | A | laptop]

location = "Hyderabad"

Q₁

Q₂

Q₃

location
cities
Hyderabad

Patna

new Delhi

Kanpur

Time
(quarters)

Q₁

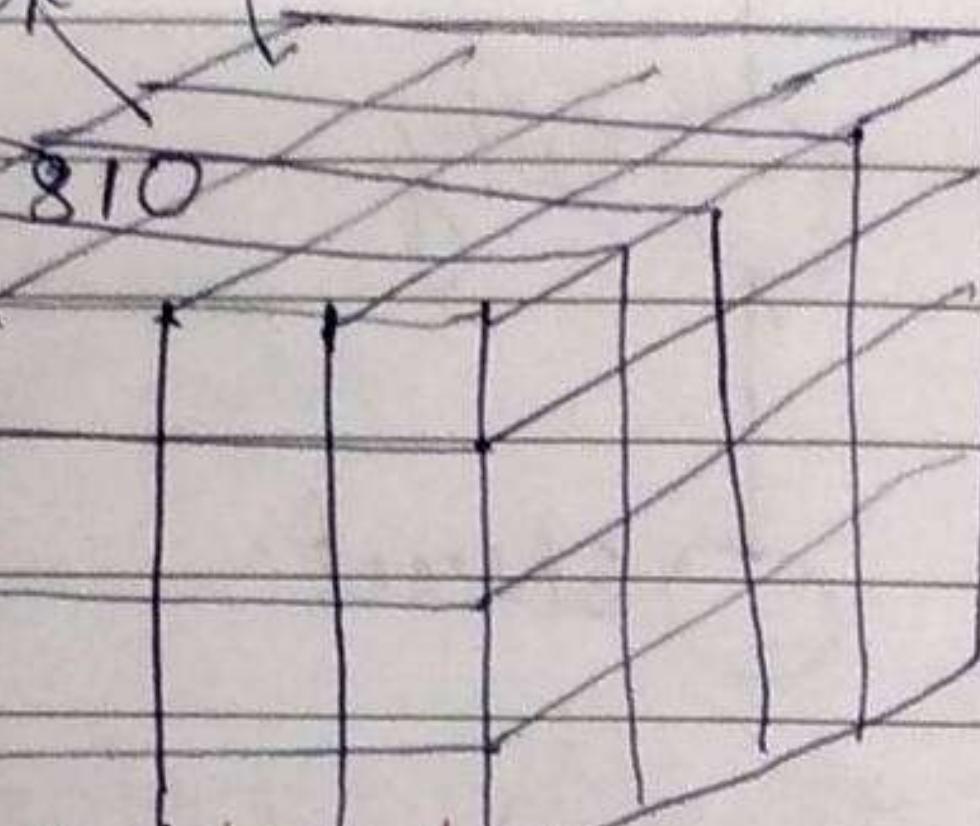
Q₂

Q₃

Q₄

Hm
En Comp.
laptop
mobile

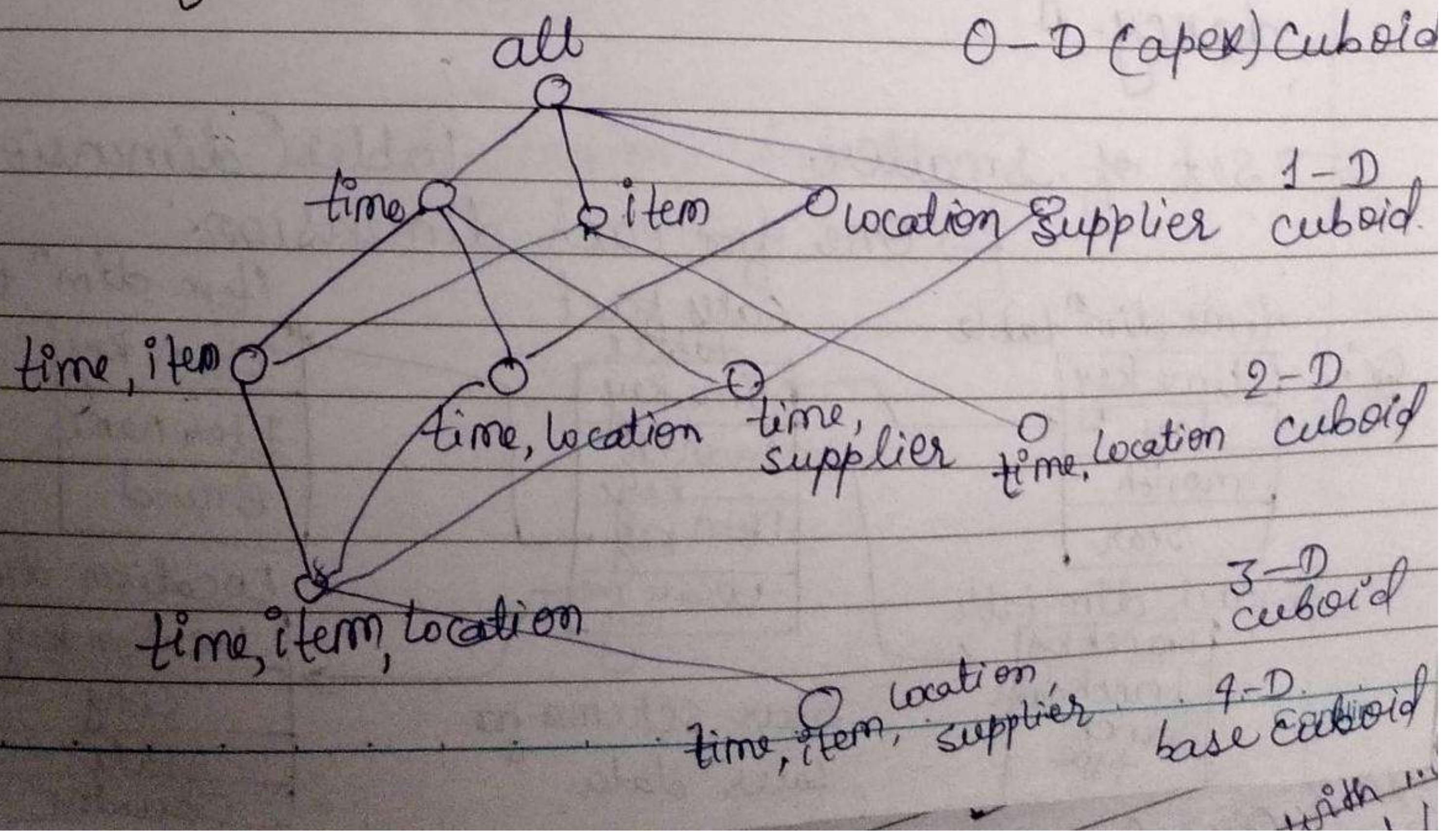
Items (types)



Add more cities.

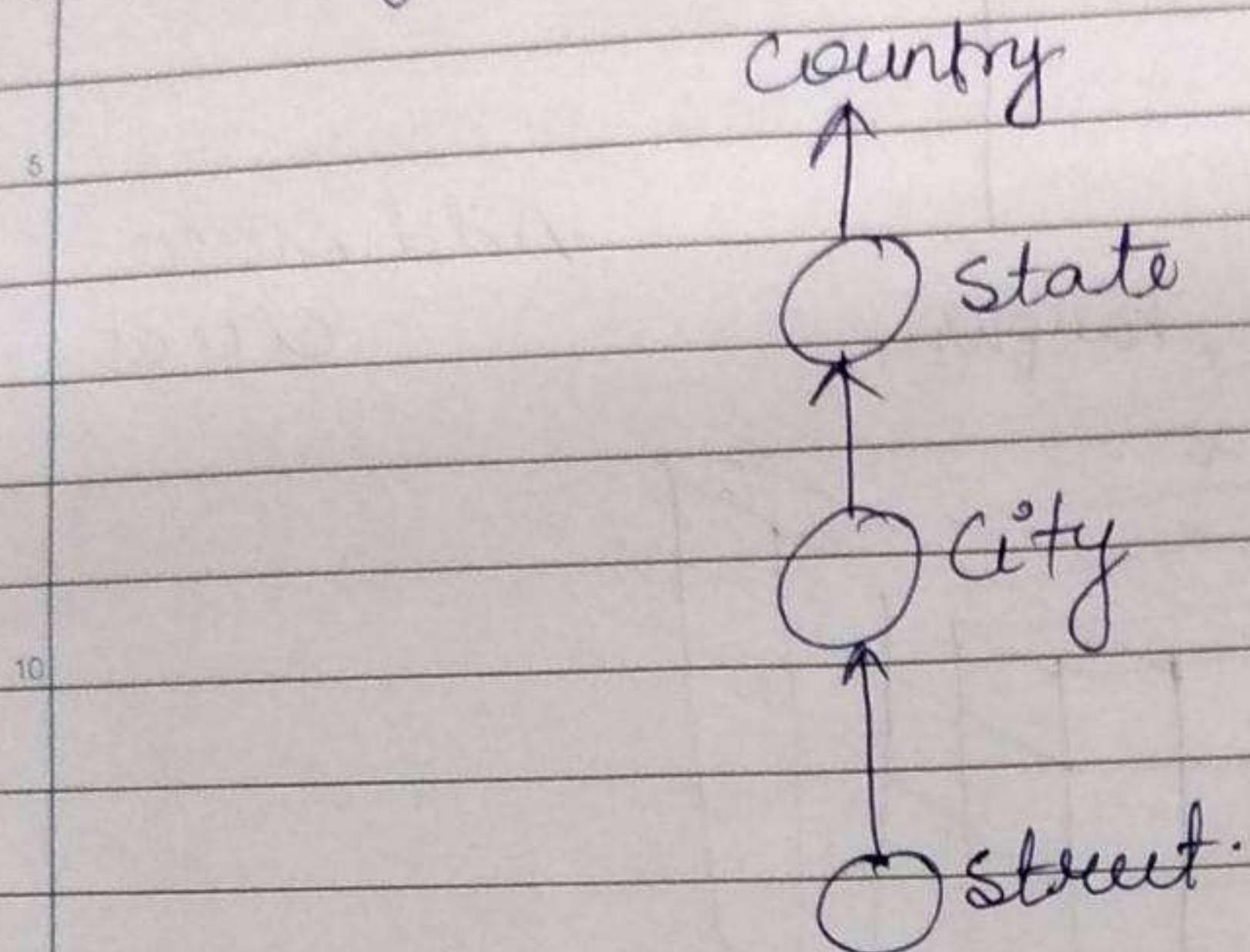
→ If we have n-D data then we can represent it as a series of (n-1)-D cubes.

→ Lattice of cuboids :-



Lattice of cuboids:

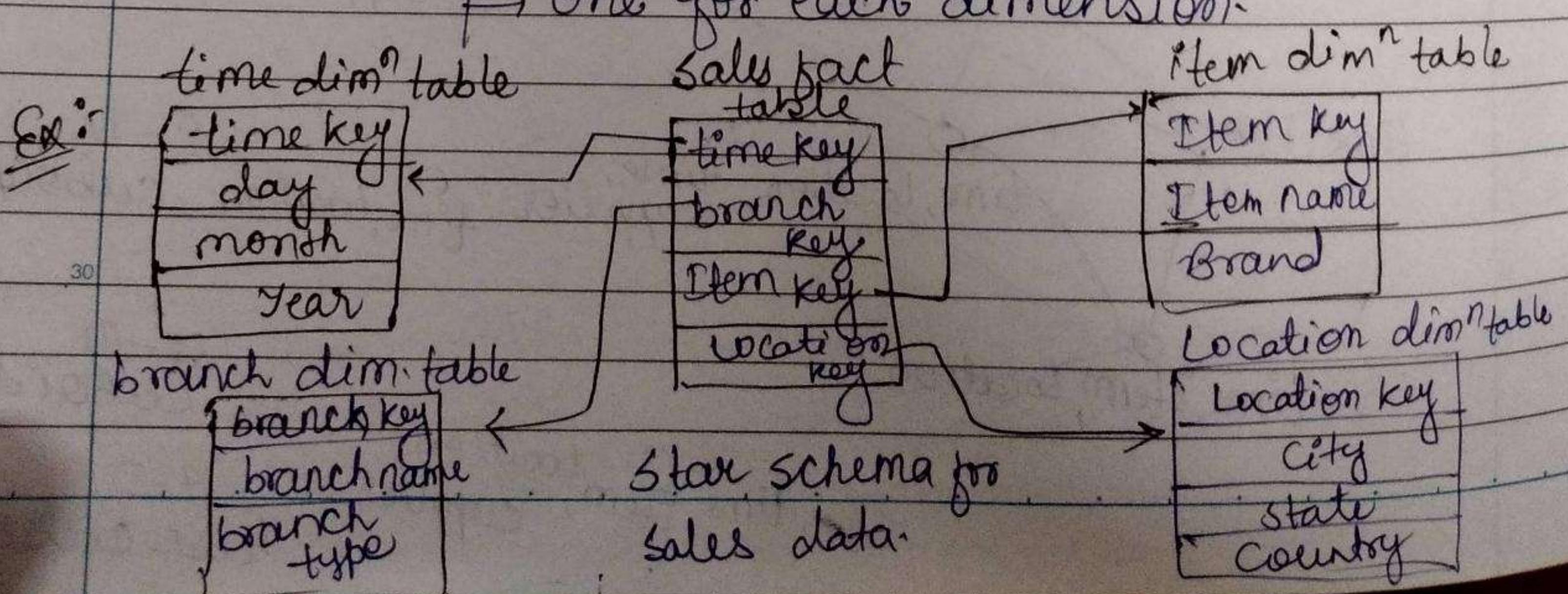
- It defines a sequence of mapping from a set of low level concepts to higher level or more general concepts.

Schemas for multidimensional databases:

- Star Schema
- Snowflake schema
- Fact constellation schema

Star schema:

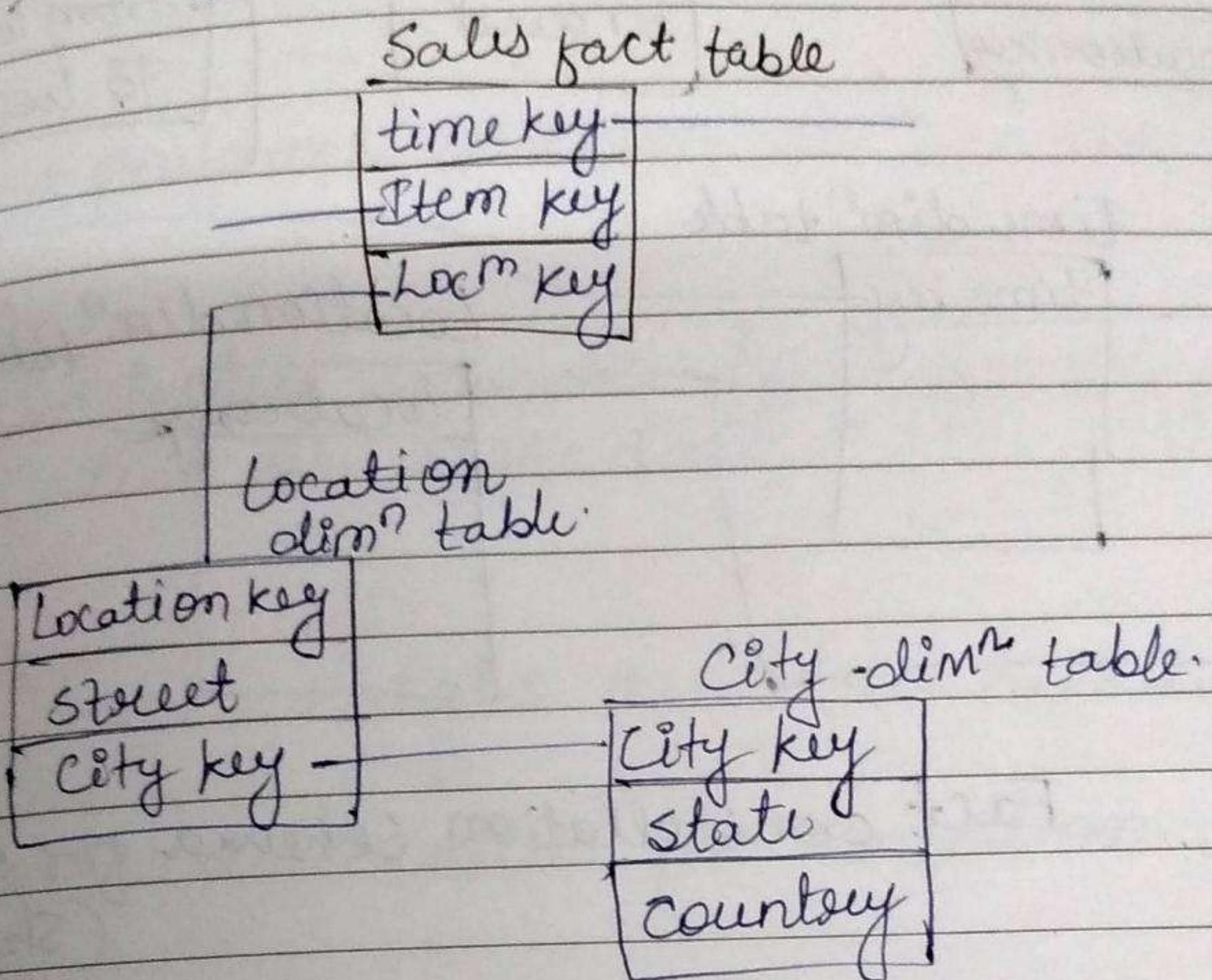
- It has large central table (fact table) containing the bulk of data with no redundancy.
- set of smaller tables (dimension tables)
 - One for each dimension.



Snowflake Schema:-

Page :
Date :

- Variant of STAR
- Some dimension table are normalized.
- further splitting of data.

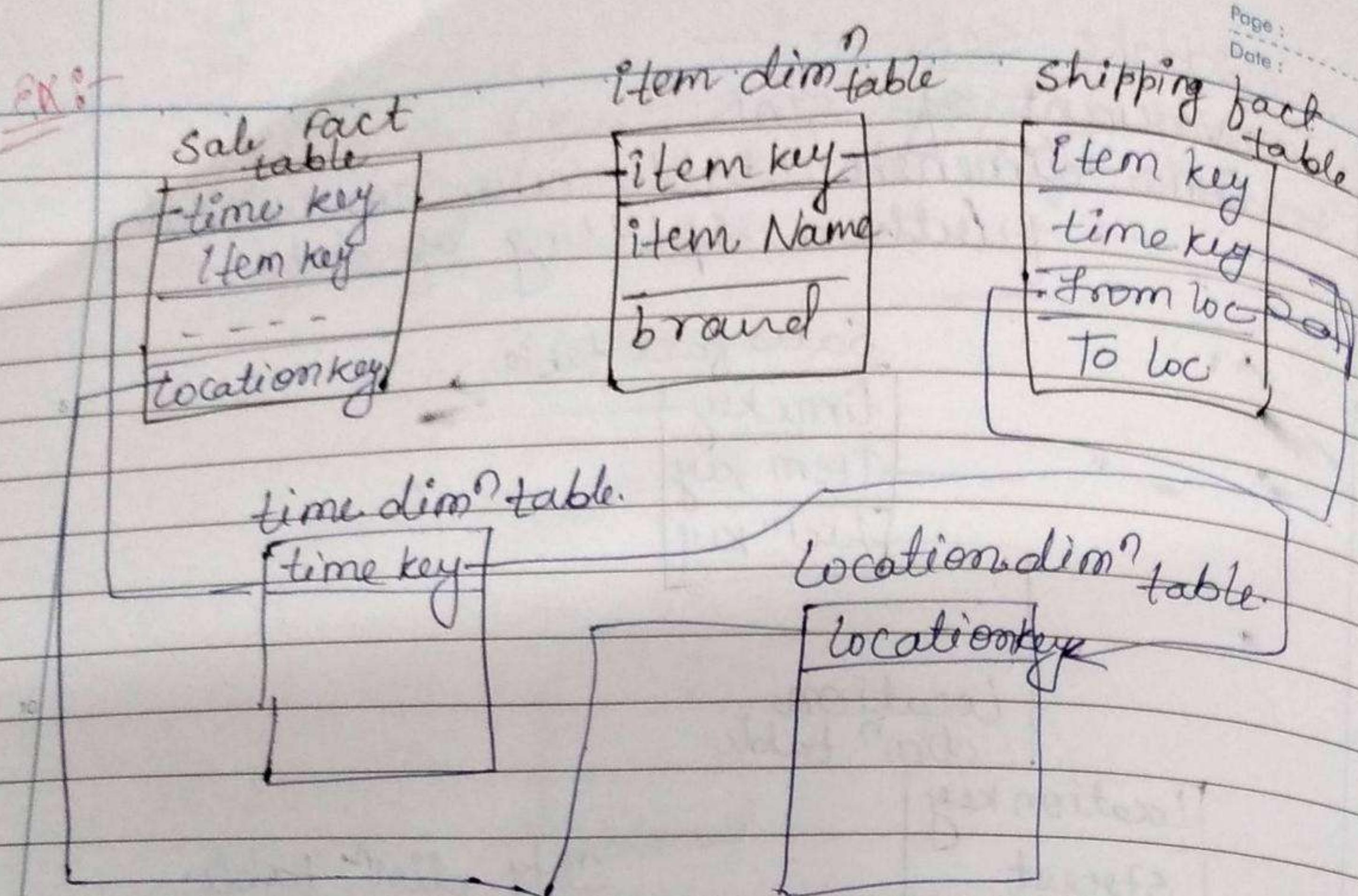


Fact constellation:-

- Galaxy schema
- Sophisticated / complex app^o.
- Has multiple fact tables that

↓
share

↓
dimension tables.



Fact constellation schema for sales
(Shipping)

→ OLAP operations :-

→ Roll-up :-

→ performs aggregation

↓
data cube

25 Climbing up the concept hierarchy
or dimⁿ reduction.

→ Drill-down :-

→ Reverse of roll-up

→ navigates from

↓
less detail

more detail

- moving down the concept hierarchy
- introduce additional dim?

→ Slice & dice :-

Slice :-

- ↳ performs selection on one dim?

Dice :-

- ↳ defines subcube by performing selection on two or more dim?

→ Pivot (rotate) :-

- ↳ It rotates the data axis in view
↓
alternative presentation.

e.g.


3 Three Level Data Warehouse Architecture

Diagram in book.

Middle tier :-

It is an OLAP server.

- ROLAP

- MOLAP

ROLAP (Relational OLAP)

- extended form of relational DBMS.

MOLAP :- (Multidimensional OLAP)

special purpose server that directly implements multidimⁿ data & operations

Top tier:-

↳ front end client layer

-> Query / reporting / Analysis tools.

Enterprise Data Warehouse

→ all information

→ about various subjects

→ entire organization

→ cross functional in scope

→ detailed & summarized data

→ Terabytes

Data Mart :-

→ subset of corporate wide data

→ specific group of users

↳ Independent

↳ Dependent

Independent :-

→ Data directly from the external information providers / other operational system.

Dependent :-

→ Data from the enterprise data warehouse.

→ Frequent Pattern

→ It appears in dataset frequently.
e.g. milk & bread.

→ Subsequence

→ Sequence that occurs frequently in the dataset.

↳ frequent sequential pattern.

PC → Pointer → Mouse Pad.

20) Market Basket Analysis :-

→ Ex: frequent itemset mining.

→ Analysis customer by buying habit.

25) → Association among diff. items that he/she places in their shopping cart.

→ Retailers to develop marketing strategies.

Data Mining & Warehousing

Association Rules

Frequent Patterns



Market Basket Analysis - is a typical example of
 ↳ Frequent itemset mining. The process
 analyzes customer buying habit by finding
 association among different items that customer
 places in their shopping basket. The discovery of
 such association can help retailers to develop
 marketing strategies by gaining insights into which
 items are purchased together by the customer.

e.g.: if a customer buys milk, how likely are
 they also to buy bread at the same time.

Association Rule:-

cart items $I_1 | I_2 | I_3 | \dots | I_n$

Customer [1 0 1 ----- 1]

If we think of the universe as a set of items available at the store, then each item has a boolean variable representing the presence/absence of that item. Each basket can then be represented by boolean vector of values assigned to these variables. The boolean vectors can then be analyzed for buying patterns that reflects items that are frequently associated purchased together. These patterns can be represented in the form of association rule.

e.g.:

customer who purchases computers are likely to buy antivirus software.

[support = 2% , confidence = 60%]

support - 2% (transactions where both are purchased)
 confidence - 60% (chances that both will be purchased)

Rules -

- * 'support' and 'confidence' are two measures of rule interestingness.
- * 2% support means 2% of all transactions under analysis show that computer and antivirus software are purchased together.
- * 60% confidence means 60% of customers who purchased a computer also bought the antivirus software.
- * Association rules are considered interesting if they satisfy both - minimum support threshold and minimum confidence threshold.
 These thresholds are provided by the users domain expert.

$\text{def } I = \{ I_1, I_2, \dots, I_n \}$ $T \rightarrow \text{transactions}$

A be a set of items

T is said to contain A if & only if
 $A \subseteq T$

an association rule $A \rightarrow B$

where $A \subset I$, $B \subset I$, and $A \cap B = \emptyset$

support $A \rightarrow B = P(A \cup B)$

confidence $A \rightarrow B = P(B|A)$

$$\text{confidence} = \frac{\text{support}(A \cup B)}{\text{support}(A)}$$

$$= \frac{\text{support-count}(A \cup B)}{\text{support-count}(A)}$$

Challenge -

suppose large ^{no} items (say 100)

$$100c_1 + 100c_2 + \dots + 100c_{100} = 2^{100} - 1$$

(very difficult to process)

- * } Closed frequent Itemset
- * } Maximal "

Closed Frequent Itemset

An itemset 'X' is closed in a dataset 'S' if \exists no proper superset 'Y' s.t. Y has the same support count as 'X' in 'S'

An itemset 'X' is Maximal if \exists no super item set 'Y' s.t. $X \subset Y$ and 'Y' is frequent in 'S'.

e.g. $C = \{ \{a_1, a_2, \dots, a_{100}\} : 1, \{a_1, a_2, \dots, a_{50}\} : 1 \}$, \rightarrow 2 transactions

min-support = 1.

Closed frequent Itemset \rightarrow

$C = \{ \{a_1, a_2, \dots, a_{50}\} : 2, \{a_1, a_2, \dots, a_{100}\} : 1 \}$

Maximal \rightarrow

$M = \{ \{a_1, a_2, \dots, a_{100}\} : 2 \}$

$\hookrightarrow a_1, \dots, a_{50}$ didn't coz. its super item set is also frequent

Multilevel Association Rules -

'X' is a variable representing customer

buys ('X', "computer") \rightarrow buys ('X', "HP Printer")

buys ('X', "Laptop") \rightarrow buys ('X', "HP Printer")

\hookrightarrow items fetched by different levels of abstraction)

Multidimension association rule (A.R)

If items/ attributes in an association rules refer only 1 dimension, then referred as 1D A.R

If the rule refers 2 or more dimensions like ~~etc~~ age, income, buys, etc., then it is multidim A.R.

e.g.: age ('x', "30--39") \wedge income ('x', "42K-48K")
buys ('x', "Sony HD LED TV")

* Apriori Algo-

frequent itemset finding using candidate generation

- R. Agrawal & R. Srikant

Used for finding frequent itemset for boolean association rules.

The name of the algorithm is based on the fact that the algo uses prior knowledge of frequent itemset properties.

It uses level wise search where k-itemsets are used to find $(k+1)$ itemsets.

First, the frequent 1-itemset is found by scanning the db, to accumulate the count of each item and collecting those items that satisfy minimum support.

The resulting set is denoted by L1.

L1 is used to find L2 that is the set of frequent 2-itemsets which is then used to find out

L3 -- and so on.

Finding of each L_k requires the complete scan of the database

Apriori property - states that if a set is frequent then its all non-empty subsets are also frequent.

If an item A is added to an itemset I then the resulting itemset $I \cup A$ cannot more frequently than I .

This property belongs to a special category of property called "Anti-Monotone" in the sense that if a set cannot pass a ~~set~~ test then all of its super set will fail in the same test as well.

L_{k-1} is used to find L_k .

It is a 2 step process consisting of a join and prune action.

Join step - to find L_k a set of candidates k -itemsets are generated by joining L_{k-1} with itself.

The set of candidates is denoted by C_k .

Apriori assumes that items within a txn are stored in the lexicographic order.

Prune Step - C_k is a superset of L_k

A scan of db is used to determine the count of each candidate in C_k ; that would result in determination of L_k .

C_k however can be large, to reduce the size of C_k the apriori property is used as follows -

Any $(k-1)$ itemsets subsets of a candidate k -itemset is not in the L_{k-1} then the candidate cannot be frequent either & so can be removed from C_k .

I.	T ₁	I ₁	I ₂	I ₅	T ₅	I ₁ , I ₃
	T ₂	I ₂	I ₄		T ₆	I ₂ , I ₃
	T ₃	I ₂	I ₃		T ₇	I ₁ , I ₃
	T ₄	I ₁	I ₂	I ₄	T ₈	I ₁ , I ₂ , I ₃ , I ₅
					T ₉	I ₁ , I ₂ , I ₃

min-support = 2

scan → DB

	I ₁	6	Compare
	I ₂	7	Candidate
	I ₃	6	Support with
	I ₄	2	min support-count
	I ₅	2	Count (2)

I ₁ , I ₂	I ₁ , I ₃
I ₁ , I ₃	I ₁ , I ₄
I ₁ , I ₄	I ₁ , I ₅
I ₁ , I ₅	I ₂ , I ₃
I ₂ , I ₃	I ₂ , I ₄
I ₂ , I ₄	I ₂ , I ₅
I ₂ , I ₅	I ₃ , I ₄
I ₃ , I ₄	I ₃ , I ₅
I ₃ , I ₅	I ₄ , I ₅

L₁ (frequent 1)C₁

I ₁ , I ₂	4	
I ₁ , I ₃	4	
I ₁ , I ₄	1	x
I ₁ , I ₅	2	
I ₂ , I ₃	4	
I ₂ , I ₄	2	
I ₂ , I ₅	2	
I ₃ , I ₄	0	x
I ₃ , I ₅	1	x
I ₄ , I ₅	0	x

Scan DB to find L₂ compare with min support to find L₂

I ₁ , I ₂	4	
I ₁ , I ₃	4	
I ₁ , I ₅	2	generate C ₂
I ₂ , I ₃	4	
I ₂ , I ₄	2	
I ₂ , I ₅	2	

I ₁ , I ₂ , I ₃	2	
I ₁ , I ₂ , I ₅	2	

L₃

I ₁ , I ₂ , I ₃	2	
I ₁ , I ₂ , I ₅	2	

END

C₂

I ₁ , I ₂ , I ₃	2	
I ₁ , I ₂ , I ₅	2	
I ₁ , I ₃ , I ₄	1	
I ₂ , I ₃ , I ₄	2	
I ₂ , I ₃ , I ₅	2	

C₃

* Generating Association Rules from frequent Itemsets

$$\text{confidence } (A \rightarrow B) = \frac{\text{support-count } (A \cup B)}{\text{support-count } (A)}$$

For every non empty subset 'S' of I output the rule
 $S \rightarrow (I - S)$ where

$$\frac{\text{support-count } (I)}{\text{support-count } (S)} \geq \text{min-confidence}$$

$$\text{min-confidence} = 70\%$$

$$I = \{I_1, I_2, I_3, I_4, I_5\}$$

$$\{I_1, I_2\} \{I_1, I_5\} \{I_2, I_5\} \{I_1\} \{I_2\} \{I_5\}$$

$$\text{confidence } I_1 \wedge I_2 \rightarrow I_5 = \frac{2}{4} \quad (\text{prev qn}) = 50\%$$

$$\text{confidence } A \rightarrow B = \frac{4}{4}$$

$$I_1, I_5 \rightarrow I_5 = \frac{2}{2} = 100\%$$

$$I_2, I_5 \rightarrow I_5 = \frac{2}{2} = 100\%$$

$$I \rightarrow I$$

FP Growth (frequent pattern) / FP Tree -

Apriori algo suffers from 2 main problems -

- ① It may need to generate huge no of candidate sets
- ② It may need to repeatedly scan the db and check the large set of candidates.

FP Growth -

The first scan of db is same as Apriori which derives the set of frequent 1-items and their support counts.

The set of frequent items are sorted in the descending order of support count. The resulting set or list is denoted by 'L'.

Firstly we create the root of the tree labeled with null.

Secondly, scan db for 2nd time and the items in each ten are processed in the 'L' order (i.e. sorted acc to the descending support count); and a branch is created for each ten.

In general, when the branch is added for a ten, the count of each node along a common prefix is incremented by 1, and the nodes for the items following the prefix are created and linked accordingly.

$$T_1 \rightarrow I_1, I_2, I_5$$

$$T_2 \rightarrow I_2, I_4$$

$$T_3 \rightarrow I_5, I_3$$

$$T_4 \rightarrow I_1, I_2, I_4$$

$$T_5 \rightarrow I_1, I_3$$

$$T_7 \rightarrow I_1, I_3$$

$$T_8 \rightarrow I_1, I_2, I_3, I_5$$

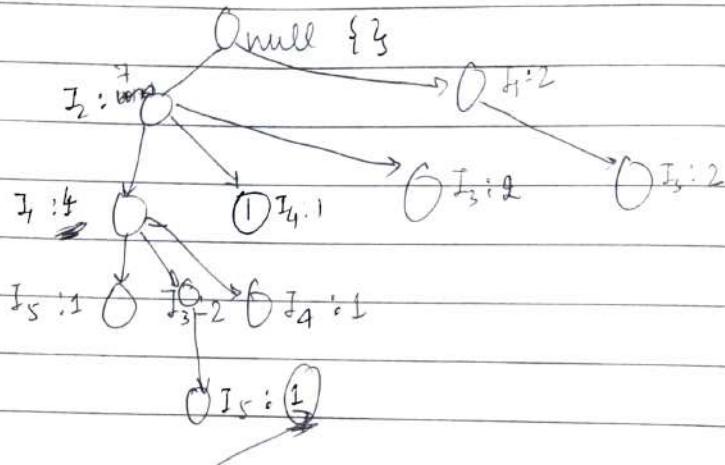
$$T_9 \rightarrow I_1, I_2, I_3$$

Find count & sort descending.

$$L = \{ \{I_2: 7\}, \{I_1: 6\}, \{I_3: 6\}, \{I_4: 2\}, \{I_5: 2\} \}$$

Now, list items in transactions @ descending order

$$\begin{array}{ll} T_1 \rightarrow I_2, I_1, I_5 & T_6 \rightarrow I_2, I_3 \\ T_2 \rightarrow I_2, I_4 & T_7 \rightarrow I_1, I_3 \\ T_3 \rightarrow I_2, I_3 & T_8 \rightarrow I_2, I_1, I_3, I_5 \\ T_4 \rightarrow I_2, I_1, I_4 & T_9 \rightarrow I_2, I_1, I_3 \\ T_5 \rightarrow I_1, I_3 & \end{array}$$



An FP Tree is mined as follows

- Start from each frequent length 1 pattern (as an initial suffix pattern) construct its conditional pattern base ; then construct its conditional FP tree and perform mining recursively in such a tree.
- The pattern growth is achieved by concatenation of the suffix pattern with the — Pattern generated from the conditional FP tree

$\therefore \text{support_count} = 2$ (minimum)

Item	Conditional Pattern Base	Conditional FP Tree	Frequent Pattern Generated
I ₅	{I ₂ : 1, {I ₁ : 1}}	{I ₂ : 2, I ₁ : 2}	{I ₂ , I ₅ : 2, I ₁ , I ₅ : 2}
I ₄	{I ₂ : 1}, {I ₂ , I ₁ : 1}	{I ₂ : 2}	{I ₂ , I ₄ : 2}
I ₃	{I ₂ : 2}, {I ₁ : 2}, {I ₂ , I ₁ : 2}	{I ₂ : 4, I ₁ : 2}	{I ₂ , I ₃ : 4}
I ₁	{I ₂ : 4}	{I ₂ : 4}	{I ₂ , I ₁ : 4}

Decision Tree

A decision tree is a flowchart-like tree structure where each internal node denotes a test on attribute and each branch represents an outcome of test & each leaf node holds a class label.

The top-most node in the tree is referred as root node.
The construction of a decision tree classifier does not require domain knowledge/ parameter settings and is suitable for exploratory knowledge discovery!

$A_1 | A_2 | \dots | A_n \rightarrow \text{attributes.}$

Attribute selection method - (ASM)

It is a procedure to determine splitting criteria.
The tree starts with a single node 'N', representing the training tuples in 'D' (dataset).

If the tuples in D are all of the same class, then node 'N' becomes a leaf node and is labeled to that class.

Otherwise the algorithm calls attribute selection method to determine the splitting criteria. It tells us which attribute to test at node 'N'.

The algo uses the same process recursively to form a decision tree for the tuples at each resulting partition, D_j .

The recursive partitioning stops when any of the following conditions are met -

① All of the tuples in partition D belongs to the same class

② There are no remaining attributes on which the tuples may be further partitioned. In this case,

majority voting is employed. This involves converting node 'N' into leaf & labeling it with most common class in 'D'.

(3) There are no tuples for a given branch

types of PSM
T
Information gain
gain Ratio
Gain Index
Gini

Information Gain - ID₃ uses information gain as its attribute selection method.

The node with highest information gain is used as the splitting attribute for node N. Expected.

info needed to classify a tuple is given by

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2 p_i$$

(p_i = probability

that an arbitrary tuple belongs to the class C_j)

p_i can be computed by $\frac{|C_j, D|}{|D|}$

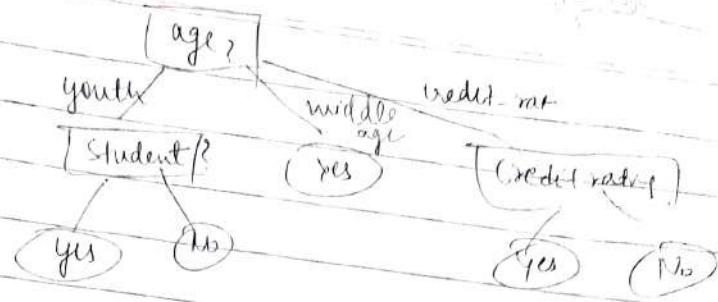
Info(D) aka Entropy of D

$$\text{Info}_n(D) = \sum_{j=1}^k \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

$|D_j|$ acts as the weight of jth partition
 $|D|$

Information Gain is defined by

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_n(D)$$



	age	income	student	credit rating	class
1)	youth	high	no	fair	buys-computer
2)	"	"	"	excellent	no
3)	middle	high	no	fair	yes
4)	senior	medium	no	fair	yes
5)	senior	low	yes	fair	yes
6)	senior	low	yes	excellent	yes
7)	middle	"	Yes	"	no
8)	youth	mediu	no	fair	yes
9)	youth	low	yes	fair	no
10)	senior	medium	yes	fair	yes
11)	youth	medium	yes	excellent	yes
12)	middle	"	No	"	yes
13)	"	high	yes	fair	yes
14)	senior	mediu	No	excellent	no

$$\text{info}(D) = \left(\frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14} \right)$$

yes no

$$\text{Info}_{\text{age}}(D) = \frac{5}{14} \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) +$$

youth yes no

middle

$$\frac{4}{14} \left(\frac{1}{4} \log_2 \frac{4}{4} + 0 \right) + \frac{5}{14} \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right)$$

senior

$$\text{Gain}(\text{Age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D)$$

$$= 0.940 - 0.694 = 0.246$$

Similarly

$$\text{Gain}(\text{Income}) = 0.029$$

$$\text{Gain}(\text{Student}) = 0.151$$

$$\text{Gain}(\text{Credit - Rating}) = 0.048$$

Gain Ratio -

The information gain measure is biased towards the test with many outcomes, i.e., it prefers to select attributes having a large no. of values.

Consider an attribute like product id / record id. A split on this attribute results in large no. of partitions. Each one contains just 1 tuple. Such a partition is sure - the info required to classify dataset 'D' based on this partition would be

$$\text{Info}_{\text{prod id}}(D) = 0$$

o The information gain by partitioning on this attribute is maximal. To avoid this

↓

C 4.5 algo (successor of ID3)

↓

Uses extension of Info Gain known as gain ratio which attempts to overcome the bias. It applies a kind of normalization to the information gain using split info

$$\text{split Info}_A(D) = -\sum_{j=1}^J \frac{D_j}{D} \log_2 \frac{D_j}{D}$$

$$\text{gain ratio}(A) = \frac{\text{Gain}(A)}{\text{split Info}(A)}$$

The attribute with splitting attribute. max gain ratio is selected as the

$$\text{Split Info (D)} = - \frac{9}{14} \log_2 \frac{9}{14} - \frac{6}{14} \log_2 \frac{6}{14} - \frac{4}{14} \log_2 \frac{4}{14}$$

9	↑	↓
high	medium	low

$$\text{Gain (Income)} = \frac{1.5789}{0.029}$$

$$\text{Gain Ratio} = \frac{0.029}{1.5789}$$

Gini Index - is used in the cart (algo)

- measures the impurity of ~~D~~ Gini(D)

$$\text{gini}(D) = 1 - \sum_{i=1}^m p_i^2$$

(where p_i is the probability that a tuple D belongs to the class C_i)

- considers the binary split of each attribute
- If a binary split on D results in D_1 and D_2 , the gini index of D is computed as

$$\text{Gini}_A(D) = \frac{|D_1|}{D} \text{gini}(D_1) + \frac{|D_2|}{D} \text{gini}(D_2)$$

* for each attribute each of the possible binary split is

the reduction in impurity that would be incurred by a binary split on a discrete or continuous valued attribute is given as -

$$\Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A(D)$$

The attribute that maximizes reduction in impurity selected as the splitting attribute

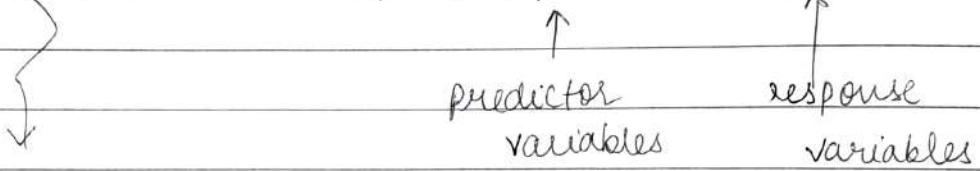
$$\text{Gini}(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2 = 0.459$$

$$\begin{aligned} \text{Gini}(D) &= \frac{8}{14} \left(1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2\right) + \frac{6}{14} \left(1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2\right) \\ &= 0.428 \end{aligned}$$

$$\Delta \text{Gini}_{\text{credit rating}} = 0.459 - 0.428$$

Regression -

Finds relation b/w independent & dependent variables.



continuous values.

Regression analyses can be used to model the relationship b/w one or more independent / predictor variables and a dependent / response variables

Linear Regression -

Straight line regression analysis involves a response variable y and a single predictor variable x . It models y as a linear fn of x .

$$y = b + wx$$

where b and w are regression coefficients.

It can also be thought of as a weight so that we can equivalently write

$$y = w_0 + w_1 x$$

$$w_1 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\sum_{i=1}^n (x_i - \bar{x})^2$$

$$w_0 = \bar{y} - w_1 \bar{x}$$

mean of x
mean of y

Q. Find the salary of employee having 10 yrs of experience.

<u>x (years of exp.)</u>	<u>Salary (y)</u>
3	30
8	57
9	64
13	72
3	36
6	43
11	59
21	90
1	20
16	83

$$y = w_0 + w_1 x \rightarrow 10$$

$$\bar{x} = \frac{554}{10} = 55.4 \quad \bar{x} = \frac{91}{10} = 9.1$$

$$w_0 = 23.6$$

$$y = 23.6 + 3.5(10)$$

Bayes Classification →

based on Bayes Thm.

All attributes are independent of each other.

$$P(H) = \frac{P(E|H) * P(H)}{P(E)}$$

where, H = hypothesis to be tested

E = evidence associated with the hypothesis

$P(A|B)$ is the conditional probability ($P(A)$ given that B has already occurred).

i.e. $P(B) > 0$

$P(H)$ is a prior probability which denotes the probability of hypotheses before presentation of any evidence

	Magazine Promotion	Watch Promotion	Life Insurance	Credit Card	Class Gender (M/F)
1	yes	No	No	No	M
2	Yes	Yes	Yes	Yes	F
3	Yes No	No	No	No	M
4	Yes	Yes	Yes	Yes	M
5	Yes	No	Yes	No	F
6	No	No	No	No	F
7	Yes	Yes	Yes	Yes	M
8	No	No	No	No	M
9	Yes	No	No	No	M
10	Yes	Yes	Yes	No	M

Find the Gender of the new instance where

magazine = Yes, watch = Yes, life ins. = No, credit card = No

gender = ? → 2-type

2 hypothesis $\rightarrow M/F$

$$P\left(\frac{\text{Gender} = M}{E}\right) = P(E | \text{Gender} = M) \cdot P(\text{Gender} = M)$$

$P(E | \text{Gender} = M)$ is a conditional probability

Computed by multiplying the conditional probab. values, for each of the evidence.

$$P(\text{Magazine} = \text{Yes} | \text{Gender} = M) = \frac{4}{6}$$

$$P(\text{Watch} = \text{Yes} | \text{Gender} = M) = \frac{2}{6}$$

$$P(\text{Life} = \text{No} | \text{Gender} = M) = \frac{1}{6}$$

$$P(\text{Credit} = \text{No} | \text{Gender} = M) = \frac{1}{6}$$

$$\Rightarrow P\left(\frac{E}{\text{Gender} = M}\right) = \left(\frac{4}{6} \times \frac{2}{6} \times \frac{1}{6} \times \frac{4}{6}\right) \times \left(\frac{6}{10}\right)$$

$$= 0.0593$$

↑ $p(\text{Gender} = M)$

No need of calculating
 $P(E)$

beccz it is common
for both

$P(E) = \text{male}$

and $P(E) = \text{female}$

$$P(\text{mag} = \text{Yes} | F) = 3/4$$

$$P(\text{watch} = \text{Yes/Female}) = 2/4$$

$$P(\text{Life} = \text{No} | \text{Female}) = 1/4$$

$$P(\text{credit} = \text{No} | \text{Female}) = 3/4$$

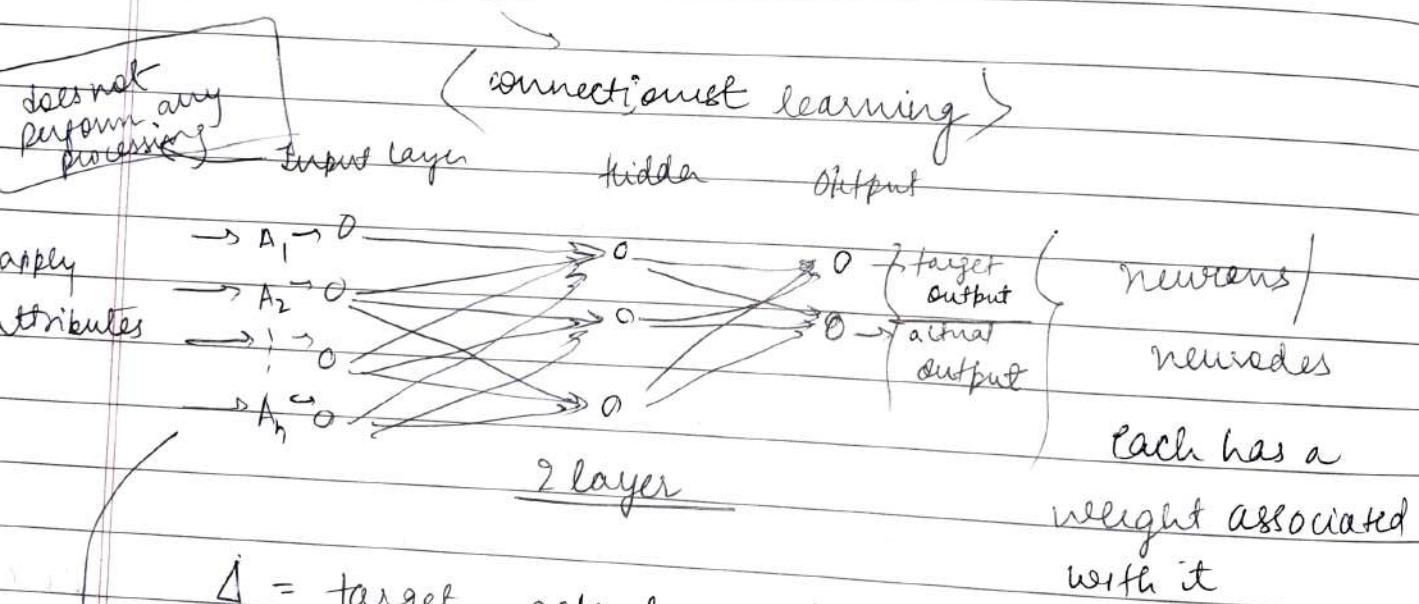
$$\frac{3}{4} \cdot \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{3}{4} \rightarrow \frac{4}{10}$$

$$\approx \underline{0.0281}$$

* This new tuple belongs to Class = Male.

$$\text{as } P(M_F) \rightarrow P(F_E)$$

* Neural Network



$$\Delta = \text{target} - \text{actual output} = \text{error.}$$

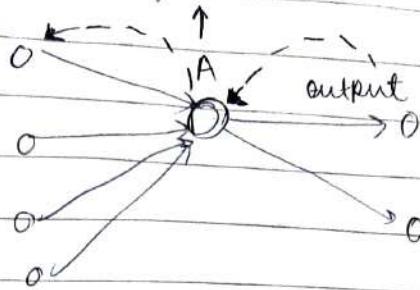
Neural Network -

based on error, we change the weights of the neural network to reduce error

→ back-propagation

Straight from neural network as no self / back edges

activation function



Back-propagation algo →

we try to find contribution of each previous layer
in the error.

And then we try to adjust weight.

A neural network is a set of connected I/O units in which each connection has a weight associated with it.

During the learning phase, the network learns by adjusting the weights associated with it, so as able to predict the correct class label of the input tuple. Neural network learning aka connectionist learning due to the connections b/w the learning.

- Back propagation algo → Performs learning in multiple layer feed forward neural learning.
It iteratively learns a set of weights for the prediction class labels of tuples.
- A multilayer feed fwd neural network consists of an input layer, one/more hidden and an output layer.
- The network is fed forward in that none of weight cycles back to an input unit or to an output of previous layer.

- For each training tuple, the weights are modified so as to minimize the mean-squared error plus the network connection & the actual target values.

The modifications are made in backward dirn i.e from output layer to each hidden layer to the first hidden layer. Hence the name back propagation.

The net input is calculated as →

particular neuron

$$I_j = \sum w_{ij} \cdot o_i + \theta_j$$

↑ ↑ bias
weight output

where w_{ij} is the weight of the connection from unit i in the prev layer to unit j .

o_i is the output of unit i from prev layer

θ_j is the bias. It serves to vary the activity of the unit.

$$o_j = \frac{1}{1 + e^{-I_j}} \quad [\text{activation fn}]$$

$I_j = \text{net input of a neuron}$

Each unit in hidden and output layer takes its net input and then applies an activation function to it. This is generally a logistic/sigmoid function.

This fn is also known as Squashing fn bcz it has large input domain in a smaller range of 0 to 1

at layer → Error $j = o_j(1 - o_j)(T_j - o_j)$

→ Error $j = o_j(1 - o_j) \sum \text{Error}_k w_{jk}$ at unit k

where w_{jk} is the weight of the from j to k

weights and biases are updated for the propagated error

Weights are updated by following eqn

$$\Delta w_{ij} = (\eta) \text{ error}_j \cdot o_j$$

$$w_{ij} = w_{ij} + \Delta w_{ij}$$

η = learning rate ($0 \leq \eta \leq 1.0$)

Bias are updated by the following eqn →

$$\Delta \theta_j = (\eta) \text{ error}_j$$

$$\theta_j = \theta_j + \Delta \theta_j$$

Cluster

Clustering is the process of grouping the data into classes or clusters, so that objects within a cluster have high similarity in comparison to one another, but are very dissimilar to objects in other clusters.

Dissimilarity are accessed based on the attribute values describing objects, often distance measures are used.

We generally organize several clustering techniques into following categories →

1) Partitioning Methods

2) Hierarchical "

3) Density based "

4) Grid " "

5) Model " "

6) Methods for high dimensional data

Constraint based clustering

	height	weight	$K=2$
x_1	64	60	
x_2	60	61	
x_3	59	70	
x_4	68	71	

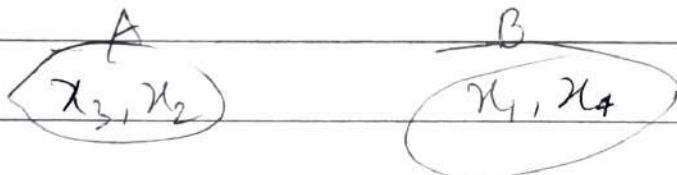
Initially (x_1, x_2) \rightarrow clusters

$$d(x_3, x_4) = \sqrt{(59-64)^2 + (70-60)^2} = 11.18$$

$$d(x_3, x_2) = 9.06$$

$$d(x_4, x_1) = 11.7$$

$$d(x_4, x_2) = 12.81$$



error = ?

$$(dis) \quad E = d(x_3, x_2)^2 + d(x_1, x_4)^2$$

computing
cluster quality

$$\text{A} \rightarrow \left(\frac{64+68}{2}, \frac{60+71}{2} \right)$$

A
 x_1, x_4

B
 x_2, x_3

$$B \rightarrow \left(\frac{60+59}{2}, \frac{61+70}{2} \right)$$

~8.32

$$\begin{array}{l|l|l|l} \text{dis}(x_1, A) = 5.85 & \text{dis}(x_2, A) = 7.5 & \text{dis}(x_3, A) = 11.18 & \text{dis}(x_4, A) = 5.85 \\ \text{dis}(x_1, B) = 7.1 & \text{dis}(x_2, B) = 4.83 & \text{dis}(x_3, B) = 11.18 & \text{dis}(x_4, B) = 10.12 \end{array}$$

$$E = d(x_1, A)^2 + d(x_4, A)^2 + d(x_1, B)^2 + d(x_3, B)^2$$

$$= 109.43$$

PAM →

↓ we only use actual data points

we take absolute dis. not square dis.
i.e.

Iteration 1) in per gn

$$E = d(x_3, x_2) + d(x_1, x_4)$$

$$= 20.76$$

Iteration 2:

The algo selects new medoids in place of existing medoids.

some other points except x_1, x_2

i.e., x_3 or x_4

x_4 as medoid in place of x_1

Cluster head →

x_2, x_4

$$d(x_1, x_2) = 4.123 \checkmark$$

$$d(x_1, x_4) = 11.7$$

$$d(x_3, x_2) = 9.06$$

$$d(x_3, x_4) = 9.05 \checkmark$$

(x_1, x_2)

(x_3, x_4)

$$E_{\text{def}} = d(x_1, x_2) + d(x_3, x_4)$$

$$= 13.17$$

If error is reduces

→ new cluster head imposed

else the previous cluster head will remain as it is.